



## Direct solvers for symmetric eigenvalue problems

Bruno Lang

published in

*Modern Methods and Algorithms of Quantum Chemistry*,  
J. Grotendorst (Ed.), John von Neumann Institute for Computing,  
Jülich, NIC Series, Vol. 1, ISBN 3-00-005618-1, pp. 203-231, 2000.

© 2000 by John von Neumann Institute for Computing

Permission to make digital or hard copies of portions of this work for personal or classroom use is granted provided that the copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise requires prior specific permission by the publisher mentioned above.

<http://www.fz-juelich.de/nic-series/>

# DIRECT SOLVERS FOR SYMMETRIC EIGENVALUE PROBLEMS

B. LANG

*Aachen University of Technology*

*Computing Center*

*Seffenter Weg 23, 52074 Aachen, Germany*

*E-mail: lang@rz.rwth-aachen.de*

This article reviews classical and recent direct methods for computing eigenvalues and eigenvectors of symmetric full or banded matrices. The ideas underlying the methods are presented, and the properties of the algorithms with respect to accuracy and performance are discussed. Finally, pointers to relevant software are given.

This article reviews classical, as well as recent state-of-the-art, direct solvers for standard and generalized symmetric eigenvalue problems. In Section 1 we explain what direct solvers for symmetric eigenvalue problems are. Section 2 describes what we may reasonably expect from an eigenvalue solver in terms of accuracy and how algorithms should be structured in order to minimize the computing time, and introduces two basic tools on which most eigensolvers are based, namely similarity transformations and deflation. For accuracy reasons, orthogonal transformations should be used whenever possible. Some simple orthogonal transformations are discussed in Section 3.

Most eigenvalue solvers work in two phases: First the matrix is reduced to tridiagonal form and then the eigenvalue problem for the tridiagonal matrix is solved. Algorithms for these two phases are discussed in Sections 4 and 5, respectively, whereas Section 6 reviews algorithms that do not rely on an initial reduction. The methods presented here also carry over to complex Hermitean matrices. For simplicity we will focus on the real case. Finally, Section 7 points to relevant software.

## 1 Setting the Stage

After introducing some notational conventions, this section recalls the symmetric standard and generalized eigenvalue problems and points out the differences between direct and iterative eigenvalue solvers.

### 1.1 Some Notation

Throughout this article, matrices are denoted by uppercase letters  $\mathbf{A}, \mathbf{B}, \dots$ , and the  $(i, j)$  entry of  $\mathbf{A}$  is referred to as  $\mathbf{A}(i, j)$ . Analogously, lowercase letters  $\mathbf{x}, \mathbf{y}, \dots$  stand for (column) vectors with entries  $\mathbf{x}(i)$ , whereas greek letters  $\alpha, \beta, \dots$  denote scalars. Unless explicitly stated otherwise, all matrices are  $n$ -by- $n$  and all vectors have length  $n$ .

$\mathbf{I}$  is the identity matrix (with entries  $\mathbf{I}(i, j) = 1$  if  $i = j$  and  $\mathbf{I}(i, j) = 0$  otherwise), and  $\mathbf{0}$  is the matrix with all entries equal to zero. For any matrix  $\mathbf{A} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{A}^\top \in \mathbb{R}^{n \times m}$  denotes the *transpose* of  $\mathbf{A}$ , i.e.,  $\mathbf{A}^\top(i, j) = \mathbf{A}(j, i)$  for all  $i, j$ . A square matrix  $\mathbf{A}$  is *symmetric* if  $\mathbf{A}^\top = \mathbf{A}$ .

## 1.2 The Symmetric Eigenvalue Problem

The *symmetric (standard) eigenvalue problem* consists of computing all or selected *eigenvalues* and associated *eigenvectors* of a symmetric matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , that is, scalars  $\lambda_i$  and vectors  $\mathbf{q}_i \neq \mathbf{0}$  satisfying

$$\mathbf{A} \cdot \mathbf{q}_i = \mathbf{q}_i \cdot \lambda_i . \quad (1)$$

The eigenvalues of  $\mathbf{A}$  are just the  $n$  roots (counting multiplicity) of its *characteristic polynomial*  $p(\lambda) = \det(\mathbf{A} - \lambda \mathbf{I})$ . (It is tempting to use this property for computing the eigenvalues: First determine the coefficients of the characteristic polynomial and then its roots. However, this method cannot be recommended because it gives highly inaccurate results.)

A pair  $(\lambda_i, \mathbf{q}_i)$  satisfying Eq. (1) is called an *eigenpair* of  $\mathbf{A}$ , and the set of all eigenvalues is called the *spectrum*,  $\text{spec}(\mathbf{A})$ . For symmetric matrices all eigenvalues are real, and there exists a complete set of  $n$  mutually orthogonal, normalized (*orthonormal*, for brevity) real eigenvectors:

$$\mathbf{q}_i^\top \mathbf{q}_j = 0 \quad \text{for } i \neq j \quad \text{and} \quad \mathbf{q}_i^\top \mathbf{q}_i = 1 \quad \text{for } i = 1, \dots, n .$$

Together with Eq. (1) this implies that  $\mathbf{A}$  has an *eigendecomposition*

$$\mathbf{A} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top , \quad (2)$$

where  $\mathbf{Q} = (\mathbf{q}_1 | \dots | \mathbf{q}_n) \in \mathbb{R}^{n \times n}$  is an *orthogonal* matrix (i.e.,  $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ ), and

$$\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_n) := \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} .$$

By convention the eigenvalues are numbered ascendingly, that is,  $\lambda_1 \leq \dots \leq \lambda_n$ .

Not all applications require computing the full eigendecomposition. Sometimes it suffices to compute only the eigenvalues, sometimes only selected eigenpairs are needed (e.g., eigenpairs for all non-negative eigenvalues or for the 100 largest eigenvalues), etc.

In the *symmetric generalized eigenvalue problem* one wants to compute eigenpairs  $(\lambda_i, \mathbf{q}_i)$  satisfying

$$\mathbf{A} \cdot \mathbf{q}_i = \mathbf{B} \cdot \mathbf{q}_i \cdot \lambda_i ,$$

where  $\mathbf{B}$  is another symmetric matrix. In most applications  $\mathbf{B}$  is also *positive definite* (i.e.,  $\mathbf{z}^\top \mathbf{B} \mathbf{z} > 0$  for all  $\mathbf{z} \neq \mathbf{0}$  or, equivalently, all eigenvalues of  $\mathbf{B}$  are positive). In our treatment of the generalized problem we will focus on this case.

## 1.3 Direct and Iterative Eigensolvers

For the solution of linear systems we can choose between direct methods like Gaussian elimination, which give the solution after a fixed number of operations, and iterative solvers like the *conjugate gradients* method, which produce a sequence of increasingly accurate approximations to the solution. Iterative solvers terminate as soon as the required precision is attained and can thus lead to considerable savings in operations (and in memory as well).

By contrast, a result from Galois theory implies that *there can be no algorithm* that computes the eigenvalues of every matrix in a finite number of operations (additions, subtractions, multiplications, divisions, roots of any order). Thus, every eigensolver must have an iterative component.

Nevertheless some of the methods are called *direct solvers*. As with linear systems, direct eigensolvers *transform* the matrix to obtain the eigensystem, whereas (purely) *iterative solvers* work with the original matrix and try to extract selected eigenvalues and eigenvectors from appropriate low-dimensional subspaces of  $\mathbb{R}^n$ . Direct solvers are the methods of choice to compute a significant portion of the eigendecomposition for small to medium-sized matrices ( $n \lesssim 5000$ , say), while iterative solvers are used when only a few (up to 100, say) eigenpairs of very large — and typically sparse — matrices are sought. In this article we discuss only the direct solvers; iterative methods are treated in another contribution.

## 2 Eigenvalue Computations: How Good, How Fast, and How?

At the beginning of this section we introduce some more notations and definitions. Then we explain why eigenvalue solvers cannot compute the exact eigenvalues and eigenvectors of a matrix, and what kind of accuracy a “good” algorithm can achieve. Issues related to minimizing the computing time on today’s high-performance computers are discussed subsequently. Finally, two basic tools for eigenvalue computations are introduced, namely similarity transformations that make the matrix in some way easier to handle, and deflation, which helps to break the problem into smaller ones whose solution is much cheaper.

### 2.1 More Notation

For contiguous portions of vectors and matrices we use the Matlab-style<sup>30</sup> colon notation:  $\mathbf{x}(i_1 : i_2)$  denotes the length- $(i_2 - i_1 + 1)$  vector consisting of  $\mathbf{x}$ ’s entries  $i_1$  through  $i_2$ , and  $\mathbf{A}(i_1 : i_2, j_1 : j_2)$  is the  $(i_2 - i_1 + 1) \times (j_2 - j_1 + 1)$  matrix containing rows  $i_1, \dots, i_2$  of  $\mathbf{A}$ ’s columns  $j_1, \dots, j_2$ . An isolated colon stands for the whole index range in the respective direction. Thus,  $\mathbf{A}(5, :)$  denotes the fifth row of  $\mathbf{A}$ , whereas  $\mathbf{A}(:, 3 : 7)$  contains columns 3 through 7 of the matrix.

The “size” of vectors and matrices is measured with *norms*. For our purposes, the *Euclidean norm*

$$\|\mathbf{x}\|_2 := \sqrt{\mathbf{x}^\top \mathbf{x}} = \sqrt{\sum_{i=1}^n \mathbf{x}(i)^2}$$

plays the dominant rôle. Its associated matrix norm is the *spectral norm*

$$\|\mathbf{A}\|_2 := \max \left\{ \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} : \mathbf{x} \neq \mathbf{0} \right\} \quad (= \max\{|\lambda_1|, |\lambda_n|\} \text{ if } \mathbf{A} \text{ is symmetric}) .$$

Since computing the spectral norm of a matrix is very expensive, often the *Frobenius*

*norm*

$$\|A\|_F := \sqrt{\sum_{i=1}^n \sum_{j=1}^n A(i,j)^2}$$

is used instead. Another important norm is the *maximum norm*

$$\|\mathbf{x}\|_\infty := \max\{|\mathbf{x}(i)| : i = 1, \dots, n\} .$$

For iterative processes the *order of convergence* indicates how fast the desired values are approached. Let  $\|\cdot\|$  denote a vector norm. A sequence of vectors  $(\mathbf{x}_k)_{k \geq 0}$  converges *linearly* to a limit  $\mathbf{x}^*$  if there is some constant  $0 < c < 1$  such that  $\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq c \cdot \|\mathbf{x}_k - \mathbf{x}^*\|$  for all  $k$ . This means that a constant number of steps is required to obtain one additional correct digit in the approximation  $\mathbf{x}_k$ . The sequence converges *quadratically* (or *cubically*) if  $\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq C \cdot \|\mathbf{x}_k - \mathbf{x}^*\|^2$  (or  $\|\mathbf{x}_{k+1} - \mathbf{x}^*\| \leq C \cdot \|\mathbf{x}_k - \mathbf{x}^*\|^3$ ) for some  $C > 0$ . With quadratic and cubic convergence, the number of correct digits is *doubled* or *tripled* in each step. Thus  $\mathbf{x}^*$  is approached very fast, once we have come somewhere close to it.

## 2.2 Accuracy Issues

Two obstacles prevent us from getting the *exact* eigenvalues and eigenvectors of a matrix on a computer.

First, in general the computation of the eigensystem involves an iterative process, which must be interrupted at some point before the correct values are reached, resulting in a so-called *truncation error*.

But even if this were not the case, the eigenvalues typically cannot be stored exactly in the finite number of memory cells that are allocated for each “real” number. The same holds for the results of most intermediate results. E.g., the IEEE standard 754<sup>2</sup> represents double-precision numbers with 64 bits, twelve of them encoding the sign and magnitude, and the remaining 52 bits holding the most significant binary digits of the number — corresponding to roughly 16 decimal places accuracy.

The best we may reasonably expect is that the computed result  $\text{fl}(\alpha \circ \beta)$  of each operation, where  $\circ$  stands for addition, subtraction, multiplication, or division, is the floating-point number that is closest to the exact result  $\alpha \circ \beta$ , that is,

$$\text{fl}(\alpha \circ \beta) = (\alpha \circ \beta) \cdot (1 + \epsilon) , \quad (3)$$

where  $\epsilon$  is some tiny quantity depending on the operation  $\circ$  and on the operands  $\alpha$  and  $\beta$ . If the arithmetic of the computer conforms to the IEEE standard (this is true for all recent workstations and personal computers) then the quantities  $\epsilon$  in Eq. (3) are guaranteed to be below some constant bound  $\varepsilon$ , the so-called *machine epsilon*. In double-precision,  $\varepsilon \approx 2.22 \cdot 10^{-16}$ .

But even if each single operation produces only a tiny relative error, a badly chosen (though mathematically correct) sequence of operations may yield an utterly wrong final result. In addition it is often very difficult to derive good estimates for the error of the computed quantity, the so-called *forward error*.

A much simpler technique for analyzing the behaviour of the algorithms was invented by Wilkinson<sup>39</sup>. Roughly speaking, his *backward error* analysis tries to put the blame for the rounding errors on the initial data. In the context of eigenvalue computations this means proving that the computed eigenvalues  $\text{fl}(\lambda_i)$  are the *exact* eigenvalues not of the original matrix  $\mathbf{A}$ , but of a (slightly) perturbed matrix  $\mathbf{A} + \Delta\mathbf{A}$ . An eigenvalue algorithm is called (*backward*) *stable* if the perturbation  $\Delta\mathbf{A}$  is very small compared to  $\mathbf{A}$ :

$$\|\Delta\mathbf{A}\|_2 = \mathcal{O}(\varepsilon) \cdot \|\mathbf{A}\|_2, \quad (4)$$

where  $\mathcal{O}(\varepsilon)$  stands for a “small multiple” of the machine epsilon that may grow with a low-degree polynomial in the matrix dimension, like  $4n\varepsilon$  or  $2n^2\varepsilon$ .

In order to derive bounds for the errors  $|\text{fl}(\lambda_i) - \lambda_i|$  ( $\lambda_i$  being the exact eigenvalues of  $\mathbf{A}$ ), the backward error analysis must be complemented by *perturbation analysis*, which investigates how much the eigenvalues can change if the matrix is perturbed. For the symmetric case, the following simple bound holds.

**Theorem<sup>18</sup>**. *Let  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  be symmetric matrices with eigenvalues  $\lambda_1 \leq \dots \leq \lambda_n$  and  $\tilde{\lambda}_1 \leq \dots \leq \tilde{\lambda}_n$ , respectively. Then*

$$|\tilde{\lambda}_i - \lambda_i| \leq \|\tilde{\mathbf{A}} - \mathbf{A}\|_2. \quad (5)$$

The accumulated changes of all eigenvalues can be estimated with the Frobenius norm.

**Theorem** (Wielandt-Hoffman)<sup>18</sup>. *Let  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  be symmetric matrices with eigenvalues  $\lambda_1 \leq \dots \leq \lambda_n$  and  $\tilde{\lambda}_1 \leq \dots \leq \tilde{\lambda}_n$ , respectively. Then*

$$\sqrt{\sum_{i=1}^n (\tilde{\lambda}_i - \lambda_i)^2} \leq \|\tilde{\mathbf{A}} - \mathbf{A}\|_F.$$

The changes of eigenvectors corresponding to *simple* eigenvalues can be bounded as follows.

**Theorem<sup>11</sup>**. *Let  $\mathbf{A}$  and  $\tilde{\mathbf{A}}$  be symmetric matrices with eigenvalues  $\lambda_1 \leq \dots \leq \lambda_n$  and  $\tilde{\lambda}_1 \leq \dots \leq \tilde{\lambda}_n$  and associated orthonormal eigenvectors  $\mathbf{q}_1, \dots, \mathbf{q}_n$  and  $\tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_n$ , respectively. Let  $\theta_i$  denote the (acute) angle between  $\tilde{\mathbf{q}}_i$  and  $\mathbf{q}_i$ . Then*

$$\frac{1}{2} \sin 2\theta_i \leq \frac{\|\tilde{\mathbf{A}} - \mathbf{A}\|_2}{\min\{|\lambda_j - \lambda_i| : j \neq i\}} \quad \text{if } \lambda_j \neq \lambda_i \text{ for all } j \neq i. \quad (6)$$

Therefore the direction of an eigenvector  $\mathbf{q}_i$  changes only little, provided that the matrix itself changes very little and that the associated eigenvalue  $\lambda_i$  is simple and well separated from the remaining eigenvalues. Note that when  $\theta_i$  is small then  $\frac{1}{2} \sin 2\theta_i \approx \sin \theta_i \approx \theta_i$ .

Eqs. (4), (5), and (6) imply that stable algorithms can compute the eigenvalues of a symmetric matrix  $\mathbf{A}$  with an error  $\mathcal{O}(\varepsilon) \cdot \|\mathbf{A}\|_2$  and the eigenvectors corresponding to well-separated eigenvalues with an error  $\mathcal{O}(\varepsilon)$ . For the large eigenvalues ( $|\lambda_i| \approx \|\mathbf{A}\|_2$ ) this means that the computed values will be correct except for a few digits at the end. For very small eigenvalues ( $|\lambda_i| \approx \mathcal{O}(\varepsilon)\|\mathbf{A}\|_2$ ), however, very few correct figures can be guaranteed. Indeed, even the magnitude and the sign of the computed eigenvalue may be wrong.

In general, the bounds given in Eqs. (5) and (6) cannot be improved much. For particular classes of matrices, however, the dependence of the eigenvalues on the matrix elements is much stronger: Then *relative* perturbations of the matrix (i.e., each matrix entry is perturbed proportionally to its magnitude) lead to relative eigenvalue perturbations. In these situations appropriate algorithms can compute almost all digits even of extremely tiny eigenvalues<sup>4</sup>.

### 2.3 Performance Issues

Among the stable methods we would like to select one that takes the least time to determine the desired eigenvalues and eigenvectors. The computing time depends on several factors, most notably on the number of operations, on the rate at which these operations can be performed, and on the potential for exploiting parallelism.

In the context of direct eigensolvers, the overall work is adequately captured by counting only the floating-point operations (**flop**: additions, subtractions, multiplications, divisions, and square roots) and ignoring everything else, like index manipulations etc.

The *performance* (or *execution rate*) of an algorithm is measured in **Mflop/s** (millions of **flop** per second). On today's machines with processors running at hundreds of MHz and main memories being almost an order of magnitude slower, the performance is mainly determined by the “data re-use factor”  $r$ , which is the number of operations performed, divided by the amount of data that are moved between the main memory and the processor (more precisely, between main memory and the caches, which are small, but fast memory buffers running almost at full processor speed). This issue is best explained with the different levels of the *BLAS* (*Basic Linear Algebra Subprograms*)<sup>13,14,28</sup>.

The BLAS define a set of subroutines for performing some simple recurring tasks in linear algebra computations. To consider just three of these routines, the (double-precision) function **DDOT** returns the scalar product  $\mathbf{x}^\top \mathbf{y}$  of two vectors. This routine does  $2n - 1$  **flop** and reads  $2n$  elements from memory (the entries of the two vectors), yielding  $r \approx 1$ . Another routine, **DGEMV**, computes a matrix–vector product with roughly  $2n^2$  **flop** and  $n^2$  accesses to memory, thus  $r \approx 2$ . And finally, the routine for computing matrix–matrix products, **DGEMM**, requires about  $4n^2$  memory accesses to do  $2n^3$  **flop**, resulting in a much higher ratio  $r \approx n/2$ . (The “level 1 BLAS” do order-of- $n$  **flop** on order-of- $n$  data, the level 2 routines do order-of- $n^2$  **flop** on order-of- $n^2$  data, and the level 3 routines do order-of- $n^3$  **flop** on order-of- $n^2$  data.) As can be seen from Table 1, a higher re-use factor can significantly improve the performance. (Note that only the level 3 routine comes anywhere close to the processor's peak performance of 266**Mflop/s**.)

Therefore the algorithms should be (re)structured in such a way that a major part of their operations can be done with level 3 routines, even if the overall number of operations is slightly increased. Algorithms with this property are often called *blocked algorithms*.

Modern high-performance computers typically feature multiple processors. To fully utilize the potential power of these machines, *parallelism* must be exploited, i.e., the work (and in most cases the data, too) must be distributed among the

Table 1. Performance of selected operations on a 266MHz PentiumII. All vectors had length  $n = 1000$ , all matrices were  $n$ -by- $n$ .

Operation	BLAS Routine	re-use factor $r$	Mflop/s
$\alpha := \mathbf{x}^\top \mathbf{y}$	DDOT (BLAS 1)	1	36.1
$\mathbf{y} := \alpha \mathbf{A} \mathbf{x} + \beta \mathbf{y}$	DGEMV (BLAS 2)	2	61.1
$\mathbf{C} := \alpha \mathbf{A} \mathbf{B} + \beta \mathbf{C}$	DGEMM (BLAS 3)	$n/2$	195.2

processors. In particular, the algorithm must contain enough *independent* operations to keep all the processors busy. In addition, a processor cannot work completely on its on, but from time to time it must *synchronize* with other processors in order to exchange information. As synchronization causes significant administrative overhead, one should strive for so-called *coarse-grained parallelism*, where synchronizations occur rarely, thus allowing each processor to do a lot of “useful” operations between them — as opposed to *fine-grained parallelism* with frequent synchronizations and just a few flop in between.

Algorithms that are based mainly on level 3 operations lend themselves in a natural way to coarse-grained parallelism. This is another argument in favor of blocking, besides the fact that these algorithms also achieve high per-node computing performance.

#### 2.4 Similarity Transformations

Almost every direct method for computing eigensystems makes use of the fact that eigenvalues are invariant under *similarity transformations*  $\mathbf{A} \mapsto \mathbf{X}^{-1} \mathbf{A} \mathbf{X} =: \tilde{\mathbf{A}}$ , where  $\mathbf{X}$  denotes an arbitrary non-singular matrix. More precisely, if  $(\lambda, \mathbf{x})$  is an eigenpair of  $\mathbf{A}$  then  $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$ , and hence

$$(\mathbf{X}^{-1} \mathbf{A} \mathbf{X}) \cdot (\mathbf{X}^{-1} \mathbf{x}) = (\mathbf{X}^{-1} \mathbf{x}) \cdot \lambda .$$

Thus  $(\lambda, \mathbf{X}^{-1} \mathbf{x})$  is an eigenpair of the transformed matrix  $\tilde{\mathbf{A}}$ .

Most eigensolvers exploit this property in the following way: First, a suitable similarity transformation  $\mathbf{A} \mapsto \mathbf{X}^{-1} \mathbf{A} \mathbf{X} =: \tilde{\mathbf{A}}$  reduces  $\mathbf{A}$  to a matrix  $\tilde{\mathbf{A}}$  that is in some way more easily handled (see Sections 4 and 5), then the desired eigenvalues  $\tilde{\lambda}_i$  and associated eigenvectors  $\tilde{\mathbf{x}}_i$  of  $\tilde{\mathbf{A}}$  are computed, and finally these are back-transformed into the eigenvalues  $\lambda_i = \tilde{\lambda}_i$  and eigenvectors  $\mathbf{x}_i = \mathbf{X} \cdot \tilde{\mathbf{x}}_i$  of the original matrix  $\mathbf{A}$ .

Whenever possible, *orthogonal* similarity transformations  $\mathbf{A} \mapsto \mathbf{Q}^\top \mathbf{A} \mathbf{Q}$  ( $\mathbf{Q}$  being an orthogonal matrix) should be used, for three reasons. First, the inverse  $\mathbf{Q}^{-1} = \mathbf{Q}^\top$  is readily available. Second, these transformations preserve symmetry: If  $\mathbf{A}$  is symmetric then  $\mathbf{Q}^\top \mathbf{A} \mathbf{Q}$  is symmetric, too. And third, orthogonal transformations are very *stable* in the sense that they induce a small backward error. It is a nice feature of the symmetric (standard) eigenvalue problem that in principle all computations can be done with orthogonal transformations.



### 2.5 Deflation

Suppose that we have found an orthogonal matrix  $Q$  such that

$$Q^T A Q = \begin{pmatrix} A_1 & 0 \\ 0 & A_2 \end{pmatrix} \quad (7)$$

with  $A_1 \in \mathbb{R}^{k \times k}$  and  $A_2 \in \mathbb{R}^{(n-k) \times (n-k)}$ . Then we have  $\text{spec}(A) = \text{spec}(A_1) \cup \text{spec}(A_2)$ , and therefore the eigenvalues of  $A$  can be obtained by computing those of the smaller matrices  $A_i$ . This reduction of the size of the matrices that must be transformed further — the so-called *deflation* — leads to considerable flop savings.

The eigenvectors of  $A$  may also be obtained from those of the smaller matrices  $A_i$ . To this end we partition  $Q$  accordingly:  $Q = (Q_1 | Q_2)$ , where  $Q_1$  consists of the first  $k$  columns in  $Q$ . Then Eq. (7) is equivalent to  $AQ_i = Q_i A_i$ ,  $i = 1, 2$ . Given an eigenpair  $(\tilde{\lambda}, \tilde{\mathbf{x}})$  of one of the smaller matrices  $A_i$ , we thus obtain

$$A \cdot (Q_i \tilde{\mathbf{x}}) = Q_i \cdot (A_i \tilde{\mathbf{x}}) = (Q_i \tilde{\mathbf{x}}) \cdot \tilde{\lambda},$$

which means that  $(\tilde{\lambda}, Q_i \tilde{\mathbf{x}})$  is an eigenpair of  $A$ .

## 3 Tools of the Trade: Basic Orthogonal Transformations

Orthogonal matrices play an important rôle in the computation of eigenvalues and eigenvectors (e.g., the eigenvector matrix  $Q$  in the eigendecomposition (2) is orthogonal). These orthogonal matrices are built up from two basic types, rotations and Householder transformations, which are introduced in this section. We also discuss techniques for improving the performance by applying these basic transformations in a blocked fashion.

### 3.1 Rotations

An  $(i, j)$  rotation in  $\mathbb{R}^n$  is given by the matrix

$$R = R(i, j, \theta) = \begin{pmatrix} 1 & & & & & & & & \\ & \ddots & & & & & & & \\ & & 1 & & & & & & \\ & & & c & & & s & & \\ & & & & 1 & & & & \\ & & & & & \ddots & & & \\ & & & & & & 1 & & \\ & & & -s & & & c & & \\ & & & & & & & 1 & \\ & & & & & & & & \ddots & \\ & & & & & & & & & 1 \end{pmatrix} \begin{matrix} \leftarrow i \\ \leftarrow j \end{matrix} \in \mathbb{R}^{n \times n}, \quad (8)$$

where  $c = \cos \theta$ ,  $s = \sin \theta$ , and all the remaining entries of  $R$  are zero. It is easy to verify that  $R$  is orthogonal.

Applying  $R$  to some column vector  $\mathbf{x} \in \mathbb{R}^n$  via  $\mathbf{x} \mapsto R^T \cdot \mathbf{x}$  corresponds to a counter-clockwise rotation by the angle  $\theta$  in the  $(i, j)$  coordinates plane. The same

is achieved for row vectors  $\mathbf{x}^\top$  via  $\mathbf{x}^\top \mapsto \mathbf{x}^\top \cdot \mathbf{R}$ . In both cases only the  $i$ th and  $j$ th entries of  $\mathbf{x}$  change.

By letting

$$c = \frac{\mathbf{x}(i)}{\sqrt{\mathbf{x}(i)^2 + \mathbf{x}(j)^2}} \quad \text{and} \quad s = \frac{-\mathbf{x}(j)}{\sqrt{\mathbf{x}(i)^2 + \mathbf{x}(j)^2}} \quad (9)$$

in the above transformations we can zero the  $j$ th entry of  $\mathbf{x}$ . (Note that only the cosine and the sine of the rotation angle  $\theta$  are present in the formulas (8) and (9); there is no need to compute the angle itself.)

Applying the rotation to a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  from the left,  $\mathbf{A} \mapsto \mathbf{R}^\top \cdot \mathbf{A}$ , affects only the  $i$ th and  $j$ th row of  $\mathbf{A}$ ,

$$\begin{aligned} \text{tmp} &:= c \cdot \mathbf{A}(i, :) + s \cdot \mathbf{A}(j, :) \\ \mathbf{A}(j, :) &:= -s \cdot \mathbf{A}(i, :) + c \cdot \mathbf{A}(j, :) \\ \mathbf{A}(i, :) &:= \text{tmp} \end{aligned}$$

whereas  $\mathbb{R}^{m \times n} \ni \mathbf{A} \mapsto \mathbf{A} \cdot \mathbf{R}$  does the same with columns  $i$  and  $j$ . In either case, the transformation takes  $4m$  multiplications and  $2m$  additions. If a symmetric matrix  $\mathbf{A}$  is transformed from both sides,  $\mathbf{A} \mapsto \mathbf{R}^\top \cdot \mathbf{A} \cdot \mathbf{R}$ , then the resulting matrix is again symmetric. Therefore the transformations must be applied only to the lower (or upper) triangle of  $\mathbf{A}$ , thereby approximately halving the overall cost to  $6n$  flop.

Rotations are very stable. It can be shown<sup>40</sup> that the backward error corresponding to a transformation  $\mathbf{A} \mapsto \mathbf{R}^\top \cdot \mathbf{A}$  is bounded by  $6\varepsilon\|\mathbf{A}\|_F$ .

There are variants of rotations — called *fast Givens rotations*<sup>35</sup> — that require only half as many multiplications as the “ordinary” rotations described above, resulting in  $4m$  flop for the transformation  $\mathbf{A} \mapsto \mathbf{R}^\top \mathbf{A}$ . These savings are achieved by an appropriate (implicit) scaling of the matrix  $\mathbf{R}$ .

In rotation-based algorithms, parallelism can be exploited in two ways. Either the work of each rotation is split among the processors (yielding a rather fine-grained parallelism), or in some situations several rotations in disjoint planes can be applied simultaneously<sup>24</sup>.

### 3.2 Householder Transformations

While rotations may be used to introduce single zeros in a vector or a matrix, Householder transformations can zero out more than one entry.

A length- $n$  *Householder transformation* is described by the matrix

$$\mathbf{H} = \mathbf{H}(\mathbf{y}) = \mathbf{I} - \mathbf{y}\tau\mathbf{y}^\top \in \mathbb{R}^{n \times n}, \quad (10)$$

where  $\mathbf{y} \in \mathbb{R}^n$  is an arbitrary vector and

$$\tau = \begin{cases} 0, & \text{if } \mathbf{y} = \mathbf{0} \\ 2/\|\mathbf{y}\|_2^2, & \text{otherwise} \end{cases} \in \mathbb{R}.$$

A short computation reveals that  $\mathbf{H}$  is orthogonal and that the transformations  $\mathbf{x} \mapsto \mathbf{H}^\top \cdot \mathbf{x}$  and  $\mathbf{x}^\top \mapsto \mathbf{x}^\top \cdot \mathbf{H}$  correspond to a reflection of  $\mathbf{x}$  at the (hyper)plane perpendicular to  $\mathbf{y}$ , thus inverting  $\mathbf{x}$ ’s component in direction  $\mathbf{y}$ . (Note that  $\mathbf{H}$  is symmetric; thus we might write  $\mathbf{H}$  instead of  $\mathbf{H}^\top$ .)

To zero out the entries  $i + 1, \dots, j$  of a given vector  $\mathbf{x} \in \mathbb{R}^n$ , we choose

$$\mathbf{y} = (0, \dots, 0, \mathbf{x}(i) \pm \|\mathbf{x}(i+1:j)\|_2, \mathbf{x}(i+1), \dots, \mathbf{x}(j), 0, \dots, 0)^\top, \quad (11)$$

where “ $\pm$ ” is chosen equal to the sign of  $\mathbf{x}(i)$  in order to avoid loss of precision through cancellation. This choice gives

$$\mathbf{H}^\top \cdot \mathbf{x} = (\mathbf{x}(1), \dots, \mathbf{x}(i-1), \mp \|\mathbf{x}(i:j)\|_2, 0, \dots, 0, \mathbf{x}(j+1), \dots, \mathbf{x}(n))^\top.$$

In particular, only the entries  $i, \dots, j$  of  $\mathbf{x}$  are affected by the transformation.

Analogously the transformations  $\mathbb{R}^{n \times m} \ni \mathbf{A} \mapsto \mathbf{H}^\top \cdot \mathbf{A}$  and  $\mathbb{R}^{m \times n} \ni \mathbf{A} \mapsto \mathbf{A} \cdot \mathbf{H}$  affect only rows (columns, resp.)  $i, \dots, j$  of  $\mathbf{A}$ : Making use of Eqs. (10) and (11), in the first case one computes

$$\mathbf{z}^\top := \tau \cdot \mathbf{y}(i:j)^\top \cdot \mathbf{A}(i:j, :) \in \mathbb{R}^{1 \times m} \quad (12)$$

$$\mathbf{A}(i:j, :) := \mathbf{A}(i:j, :) - \mathbf{y}(i:j) \cdot \mathbf{z}^\top. \quad (13)$$

The matrix–vector product (12) and the rank-1-update (13) each require approximately  $2\ell m$  operations ( $\ell = j - i + 1$  is the “active length” of the transformation), for a total of  $4\ell m$  flop. (Note that one *never* computes  $\mathbf{H}^\top \cdot \mathbf{A}$  as a matrix–matrix product, which would require approximately  $2\ell^2 m$  flop.)

As with rotations, symmetry can be exploited to save operations. The two-sided transformation  $\mathbb{R}^{n \times n} \ni \mathbf{A} \mapsto \mathbf{H}^\top \cdot \mathbf{A} \cdot \mathbf{H}$  of a symmetric matrix  $\mathbf{A}$  is performed as follows:

$$\mathbf{z} := \mathbf{A} \cdot \mathbf{y} \cdot \tau \quad (14)$$

$$\begin{aligned} \mathbf{v} &:= \mathbf{z} - \mathbf{y} \cdot \frac{\tau(\mathbf{y}^\top \mathbf{z})}{2} \\ \mathbf{A} &:= \mathbf{A} - \mathbf{y} \cdot \mathbf{v}^\top - \mathbf{v} \cdot \mathbf{y}^\top. \end{aligned} \quad (15)$$

Due to the symmetry, only one triangle of  $\mathbf{A}$  must be updated in the symmetric rank-2-update (15), provided that only one triangle is used in the symmetric matrix–vector product (14). Again, the two-sided transformation of  $\mathbf{A}$  can be effected at the same cost as a one-sided transformation of a non-symmetric matrix, that is,  $4n^2$  flop.

Householder transformations are very stable, too<sup>20</sup>. The backward error corresponding to a transformation  $\mathbf{A} \mapsto \mathbf{H}^\top \cdot \mathbf{A}$  is bounded by  $c n \varepsilon \|\mathbf{A}\|_F$ , where  $c$  is some small constant.

With Householder transformations, parallelism is exploited by working on a distributed matrix  $\mathbf{A}$  and doing each of the steps (14) and (15) in parallel.

### 3.3 Blocked Householder Transformations

Applying single rotations or Householder transformations does not allow using any level 3 BLAS. The situation changes if a *sequence* of transformations must be applied to the same matrix.

Bischof and Van Loan<sup>6</sup> discovered that the product  $\mathbf{Q} = \mathbf{H}_1 \dots \mathbf{H}_k$  of  $k$  length- $n$  Householder transformations  $\mathbf{H}_j = \mathbf{I} - \mathbf{y}_j \tau_j \mathbf{y}_j^\top$  can be written in the form

$$\mathbf{Q} = \mathbf{I} - \mathbf{W} \mathbf{Y}^\top \quad (\text{WY representation}), \quad (16)$$

where  $W$  and  $Y$  are suitable  $n$ -by- $k$  matrices.

For  $k = 1$  the representation (16) is obviously valid with  $W = \mathbf{y}_1 \tau_1$  and  $Y = \mathbf{y}_1$ . To proceed from  $k$  to  $k + 1$ , we assume that the matrices  $W$  and  $Y$  in (16) are known and that  $\tilde{Q} = H_1 \cdot \dots \cdot H_{k+1} = Q \cdot H_{k+1}$ . Then a short calculation reveals that  $\tilde{Q} = I - \tilde{W}\tilde{Y}^\top$  with the  $n$ -by- $(k + 1)$  matrices

$$\tilde{W} = (W \mid Q\mathbf{y}_{k+1}\tau_{k+1}) \quad \text{and} \quad \tilde{Y} = (Y \mid \mathbf{y}_{k+1}) .$$

Thus each additional transformation  $H_j$  requires only appending one new column to  $W$  and  $Y$ . Note that in particular the columns of  $Y$  are just the vectors  $\mathbf{y}_j$  defining the transformations  $H_j$ .

Given the representation (16), applying the  $k$  transformations to a matrix  $A$  amounts to two matrix-matrix products,

$$H_k^\top \cdot \dots \cdot H_1^\top \cdot A = Q^\top \cdot A = A - Y \cdot (W^\top \cdot A) ,$$

which performs significantly better than applying the single Householder transformations and requires only marginally more operations.

Later, Schreiber and Van Loan<sup>36</sup> refined the  $WY$  representation to

$$Q = I - YTY^\top \quad (\text{compact } WY \text{ representation}), \quad (17)$$

where  $Y$  is again  $n$ -by- $k$  and  $T$  is an upper triangular  $k$ -by- $k$  matrix. Here, the case  $k = 1$  is covered by setting  $Y = \mathbf{y}_1$  and  $T = \tau_1$ , and the step from  $k$  to  $k + 1$  is done by letting

$$\tilde{Y} = (Y \mid \mathbf{y}_{k+1}) \quad \text{and} \quad \tilde{T} = \left( \begin{array}{c|c} T & -\tau_{k+1} T \mathbf{y}_{k+1} \\ \hline \mathbf{0}^\top & \tau_{k+1} \end{array} \right) .$$

The compact  $WY$  representation costs significantly less additional storage than the original  $WY$  representation does ( $T$  instead of  $W$ ), but applying the transformations in the compact representation requires a third matrix-matrix product.

Under favorable conditions, applying a sequence of rotations to a matrix can also be organized in such a way that most of the work is done in matrix-matrix products<sup>27</sup>.

## 4 Phase I: Reduction to Tridiagonal Form

Most eigenvalue solvers first reduce the symmetric matrix  $A$  to a symmetric tridiagonal matrix  $T$  (i.e.,  $T(i, j) = 0$  whenever  $|i - j| > 1$ ) in order to make the ensuing iterative process simpler and cheaper. This section first presents the standard algorithm for tridiagonalizing a full matrix, as well as a blocked variant of the algorithm. Then the back-transformation of the eigenvectors is discussed, which is a final step after computing the eigenvalues and eigenvectors of the tridiagonal matrix  $T$  (cf. Section 5). Specialized reduction methods for banded matrices and generalized eigenvalue problems are described at the end of the section.

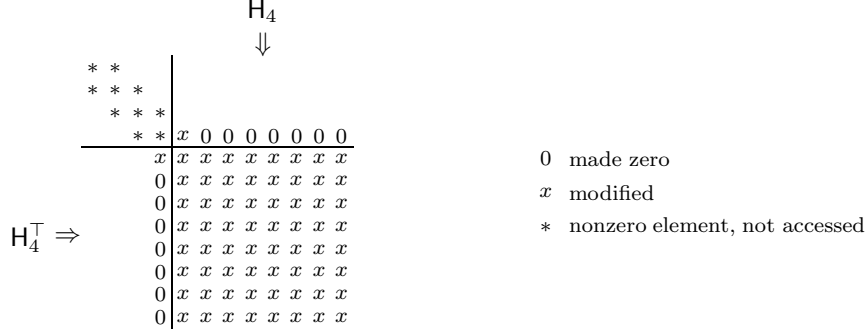


Figure 1. The fourth step in the reduction of a 12-by-12 matrix to tridiagonal form.

#### 4.1 Householder Tridiagonalization

The standard algorithm for reducing an  $n$ -by- $n$  symmetric full matrix to tridiagonal form,  $A \mapsto Q_1^\top A Q_1 = T$ , proceeds in  $n - 2$  steps:

$$A =: A_0 \mapsto H_1^\top A_0 H_1 =: A_1 \mapsto H_2^\top A_1 H_2 =: A_2 \mapsto \dots \mapsto H_{n-2}^\top A_{n-3} H_{n-2} =: A_{n-2} = T, \quad (18)$$

the  $k$ th step transforming column and row  $k$  to the desired shape with a suitable Householder transformation  $H_k$ , cf. Figure 1 and Algorithm 1.

---

**Algorithm 1** Householder reduction to tridiagonal form.

---

**for**  $k = 1$  **to**  $n - 2$   
  determine a Householder transformation  $H_k = I - \mathbf{y}_k \tau_k \mathbf{y}_k^\top$  that reduces  $A(k+1 : n, k)$  to the form  $(x, 0, \dots, 0)^\top$   
   $\mathbf{z}_k := A \mathbf{y}_k \tau_k$   
   $\mathbf{v}_k := \mathbf{z}_k - \mathbf{y}_k \cdot (\tau_k (\mathbf{y}_k^\top \mathbf{z}_k) / 2)$   
   $A := A - \mathbf{y}_k \cdot \mathbf{v}_k^\top - \mathbf{v}_k \cdot \mathbf{y}_k^\top$

---

Note that only the submatrix  $A(k+1 : n, k+1 : n)$  (and, due to symmetry, only one triangle of this matrix) is modified in the last line of the algorithm. Thus the reduction requires approximately  $\frac{4}{3}n^3$  flop, which are done mainly within the level 2 BLAS.

Before Householder developed the method just described, the tridiagonalization had been done with a rotation-based algorithm due to Givens<sup>40</sup>. Householder's approach is superior for full matrices because it requires fewer operations, whereas Givens' algorithm can make better use of sparsity.

#### 4.2 Blocked Householder Tridiagonalization

It is not necessary to build all the intermediate matrices  $A_k$  from Eq. (18) explicitly<sup>15</sup>. Instead, they can be represented in the so-called *factored form*

$$A_s = A_0 - \mathbf{y}_1 \mathbf{v}_1^\top - \mathbf{v}_1 \mathbf{y}_1^\top - \dots - \mathbf{y}_s \mathbf{v}_s^\top - \mathbf{v}_s \mathbf{y}_s^\top = A_0 - Y_s V_s^\top - V_s Y_s^\top, \quad (19)$$

where the vectors  $\mathbf{y}_k$  and  $\mathbf{v}_k$  are defined as in Algorithm 1, and  $Y_s = (\mathbf{y}_1 | \dots | \mathbf{y}_s)$  and  $V_s = (\mathbf{v}_1 | \dots | \mathbf{v}_s)$  are  $n$ -by- $s$  matrices. If the matrices  $A_k$  are built only every  $n_b$ th step then we arrive at Algorithm 2 (lines 4 through 6 and line 8 of the algorithm reflect the fact that  $\mathbf{v}_k$  and  $\mathbf{z}_k$  are computed from  $A_{k-1}$  and not from  $A_{s-1}$ , which is currently held in  $A$ ).

---

**Algorithm 2** Blocked Householder tridiagonalization.

---

```

for  $s = 1$  to  $n - 2$  step  $n_b$ 
     $Y := ()$ ,  $V := ()$  { matrices with 0 columns }
    for  $k = s$  to  $\min\{s + n_b - 1, n - 2\}$ 
        if  $k > s$ 
            compute the  $k$ th column of  $A_{k-1}$  according to Eq. (19):
             $A(:, k) := A(:, k) - YV^\top(:, k) - VY^\top(:, k)$ 
            determine the Householder transformation  $H_k$  as in Algorithm 1
             $\mathbf{z}_k := A\mathbf{y}_{k\tau_k} - Y(V^\top \mathbf{y}_{k\tau_k}) - V(Y^\top \mathbf{y}_{k\tau_k})$ 
            compute  $\mathbf{v}_k$  as in Algorithm 1
             $Y := (Y | \mathbf{y}_k)$ ,  $V := (V | \mathbf{v}_k)$ 
         $A := A - YV^\top - VY^\top$  { after  $n_b$  steps rebuild  $A$  according to Eq. (19) }

```

---

Note that the rank-2-updates (last line in Algorithm 1) have been replaced with matrix-matrix products, whereas the symmetric matrix-vector products  $A \cdot \mathbf{y}$  in the computation of  $\mathbf{z}_k$  persist. Therefore, the blocked algorithm does roughly one half of its  $\frac{4}{3}n^3$  operations with level-3 BLAS while the remaining operations are still confined to the level-2 BLAS.

The portion of matrix-matrix operations can be further increased if the reduction to tridiagonal form is done in two phases<sup>5</sup>. First the matrix is reduced to banded form (almost completely with level 3 BLAS), and then the banded matrix is tridiagonalized (no level 3 BLAS, but significantly lower flop count than for the first phase). This approach typically outperforms the direct tridiagonalization if *no eigenvectors* are required.

#### 4.3 Back-Transformation of the Eigenvectors

After the reduction of  $A$  to tridiagonal form,  $A \mapsto Q_1^\top A Q_1 = T$ , all or selected eigenpairs  $(\lambda_i, \mathbf{v}_i)$  of the tridiagonal matrix  $T$  are computed. While the eigenvalues  $\lambda_i$  of  $T$  are also eigenvalues of  $A$ , the associated eigenvectors must be back-transformed in order to obtain  $A$ 's eigenvectors via  $\mathbf{v}_i \mapsto Q_1 \cdot \mathbf{v}_i = \mathbf{q}_i$ , cf. Section 2.4.

Let  $H_k = I - \mathbf{y}_k \tau_k \mathbf{y}_k^\top$ ,  $k = 1, \dots, n - 2$ , be the Householder transformations that were used in the reduction  $A \mapsto T$ , i.e.,  $Q_1 = H_1 \cdot \dots \cdot H_{n-2}$ , and let  $V =$

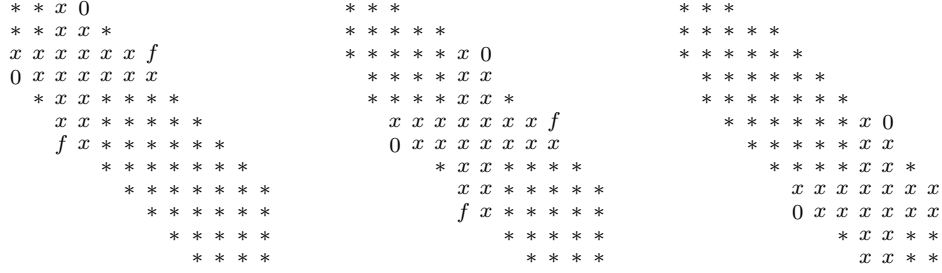


Figure 2. Rotations for zeroing the (4,1) entry and for chasing the intermediate fill-in elements (denoted by  $f$ ) in Schwarz' algorithm.

$(\mathbf{v}_1 \mid \dots \mid \mathbf{v}_m)$  denote those eigenvectors of  $\mathbf{T}$  that were computed and must be back-transformed. Then the corresponding eigenvectors of  $\mathbf{A}$  are obtained as

$$(\mathbf{q}_1 \mid \dots \mid \mathbf{q}_m) =: \mathbf{Q} = \mathbf{Q}_1 \cdot \mathbf{V} = \mathbf{H}_1 \cdot \dots \cdot \mathbf{H}_{n-2} \cdot \mathbf{V}.$$

In contrast to the tridiagonal reduction, almost all of the roughly  $2n^2m$  operations can be done with matrix-matrix products if we resort to the (compact or original) WY representation for applying  $n_b > 1$  transformations at a time, cf. Algorithm 3.

---

**Algorithm 3** Blocked back-transformation of selected eigenvectors.

---

$\mathbf{Q} := \mathbf{V}$

**for**  $s = n - 2$  **to** 1 **step**  $-n_b$

determine the compact WY representation for the next  $n'_b = \min\{n_b, s\}$

transformations:  $\mathbf{H}_{s-n'_b+1} \cdot \dots \cdot \mathbf{H}_s = \mathbf{I} - \mathbf{Y}\mathbf{T}\mathbf{Y}^\top$

apply these transformations via  $\mathbf{Q} := \mathbf{Q} - \mathbf{Y}\mathbf{T}\mathbf{Y}^\top \mathbf{Q}$

---

Like any algorithm involving mainly products of large matrices, the back-transformation can be easily and efficiently parallelized.

#### 4.4 Reduction of Banded Matrices

A symmetric matrix  $\mathbf{A}$  is *banded* with *semibandwidth*  $b$  if  $\mathbf{A}(i, j) = 0$  whenever  $|i - j| > b$ . For narrow-banded matrices ( $b \ll n$ ), Algorithms 1 and 2 are not optimal because they completely destroy the sparsity.

Such matrices are reduced with a rotation-based algorithm by Schwarz<sup>37</sup>. To understand how this algorithm works we consider a 12-by-12 matrix with semibandwidth  $b = 3$ , see Figure 2.

First, the outmost entry (4,1) in the first column is made zero with a rotation  $\mathbf{A} \mapsto \mathbf{R}^\top \mathbf{A} \mathbf{R}$  in the (3,4) plane (left picture in Figure 2). This creates a new *fill-in* entry  $f$  at position (7,3), just outside the band. Then a second rotation in the (6,7) plane is used to remove the fill-in entry (center picture), only to have another fill element appear at position (10,6), which in turn is removed by the next rotation,

and so on. Each rotation “chases” the fill element  $b$  positions down the band until the end of the matrix is reached and no further fill-in is created (right picture).

Then we can zero the next entry  $(3, 1)$  in the first column with a  $(2, 3)$  rotation, again followed by a sequence of rotations for chasing the fill-in. When the first column of  $A$  is reduced to tridiagonal form we repeat the procedure for the second column, and so on. The whole method is summarized in Algorithm 4.

---

**Algorithm 4** Schwarz’ algorithm for tridiagonalizing banded matrices.

---

```

for  $j = 1$  to  $n - 2$                                      { proceed by columns }
  for  $d = \min\{b, n - j\}$  to  $2$  step  $-1$     { zero the entry in the  $d$ th subdiagonal }
    make  $A(j + d, j)$  zero with a suitable rotation in the  $(j + d - 1, j + d)$  plane
    while the most recent rotation created a fill-in entry at some position  $(k, \ell)$ 
      make this entry zero by a suitable rotation in the  $(k - 1, k)$  plane

```

---

If  $A$ ’s eigenvectors are needed, too, then in the banded case the orthogonal matrix  $Q_1$  (cf. Section 4.3), which is the product of all rotations that are used for the reduction, is built explicitly *during* the reduction. This is achieved by applying each plane- $(k - 1, k)$  rotation  $R$  not only to the banded matrix,  $A \mapsto R^T A R$ , but also to the columns  $k - 1$  and  $k$  of an  $n$ -by- $n$  matrix  $Q_1$  via  $Q_1 \mapsto Q_1 R$ , where  $Q_1$  has been initialized as the identity matrix.

The reduction of  $A$  requires  $6bn^2$  flop. If the eigenvectors of  $A$  must be computed then the costs for accumulating  $Q_1$ ,  $3n^3$  flop, by far dominate the reduction costs. In this method some parallelism can be exploited by applying several rotations simultaneously<sup>24</sup>.

A more recent reduction algorithm is based on Householder transformations<sup>32</sup>, each transformation affecting just  $b$  rows and columns of  $A$ . This method allows coarser-grained parallelism than Schwarz’ algorithm<sup>26</sup>, and the accumulation of  $Q_1$  can be done in a blocked fashion<sup>5</sup>.

#### 4.5 Reduction of the Generalized Eigenvalue Problem

The symmetric generalized eigenvalue problem  $A\mathbf{q}_i = B\mathbf{q}_i\lambda_i$  with a symmetric *positive definite* matrix  $B$  can be transformed into a symmetric standard eigenvalue problem as follows. Let

$$B = LL^T, \quad (20)$$

where  $L$  is a lower triangular matrix with positive diagonal entries, be the *Cholesky decomposition* of  $B$ . (Algorithms for computing the *Cholesky factor*  $L$  are given below.) Then the condition  $A\mathbf{q}_i = B\mathbf{q}_i\lambda_i$  is equivalent to

$$(L^{-1}AL^{-T}) \cdot (L^T \mathbf{q}_i) = (L^{-1}BL^{-T}) \cdot (L^T \mathbf{q}_i) \cdot \lambda_i = (L^T \mathbf{q}_i) \cdot \lambda_i,$$

where  $L^{-T}$  is a shorthand for  $(L^T)^{-1} (= (L^{-1})^T)$ . Therefore, an eigenpair  $(\lambda_i, \mathbf{q}_i)$  of the generalized eigenvalue problem corresponds to the eigenpair  $(\lambda_i, L^T \mathbf{q}_i)$  of the symmetric standard eigenvalue problem for the matrix  $M = L^{-1}AL^{-T}$ . Thus we arrive at Algorithm 5, which takes approximately  $12n^3$  flop for solving the



generalized eigenvalue problem. (Note that the explicit calculation of  $L^{-1}$  and  $L^{-\top}$  is avoided by solving triangular systems with multiple right-hand sides, which is implemented in the level 3 BLAS routine DTRSM.)

---

**Algorithm 5** Generalized eigenvalue problem via Cholesky decomposition of  $B$ .

---

compute the Cholesky decomposition  $B = LL^\top$  { Algorithm 7 }  
compute  $M = L^{-1}AL^{-\top}$  by solving two triangular systems with multiple  
right-hand sides,  $XL^\top = A$  (for  $X$ ), and  $LM = X$  (for  $M$ )  
use blocked Householder tridiagonalization (Algorithm 2) and an algorithm from  
Section 5 to compute the eigendecomposition  $M = \tilde{Q}\Lambda\tilde{Q}^\top$   
compute the eigenvectors  $q_i = L^{-\top}\tilde{q}_i$  of the generalized problem by solving the  
triangular system  $L^\top Q = \tilde{Q}$  for  $Q$

---

In general, the matrix  $M$  will be full even if  $A$  and  $B$  (and therefore  $B$ 's Cholesky factor  $L$ , too) are *banded*. An algorithm by Crawford<sup>9</sup> avoids building  $M$  explicitly by interleaving its computation with the ensuing reduction to tridiagonal form. This leads to considerable flop and memory savings.

If  $B$  is ill-conditioned (i.e., its eigenvalues vary over many orders of magnitude) then severe loss of accuracy may happen because the backward error of Algorithm 5 grows with  $\|B^{-1}\|_2$ . In this case the Cholesky decomposition should be replaced with the *eigendecomposition*  $B = S\Delta^2S^\top = (S\Delta)(S\Delta)^\top$ , where  $S$  is orthogonal and  $\Delta$  contains the roots of  $B$ 's (positive!) eigenvalues. This leads to  $M = \Delta^{-1}S^\top AS\Delta^{-1}$ . In practice, the use of orthogonal matrices gives better results, albeit at higher cost.

The generalized eigenvalue problem with a matrix  $B$  that is not positive definite requires completely different techniques<sup>18,34</sup>.

Let us finally give two algorithms for computing the Cholesky decomposition. Eq. (20) is equivalent to

$$\begin{aligned} B(1,1) &= L(1,1) \cdot L(1,1) , \\ B(2:n,1) &= L(2:n,1) \cdot L(1,1) , \text{ and} \\ B(2:n,2:n) &= L(2:n,2:n) + L(2:n,1) \cdot L(2:n,1)^\top , \end{aligned}$$

i.e.,  $L(1,1) = \sqrt{B(1,1)}$ ,  $L(2:n,1) = B(2:n,1)/L(1,1)$ , and  $L(2:n,2:n)$  is the Cholesky factor of  $B(2:n,2:n) - L(2:n,1) \cdot L(2:n,1)^\top$ . Resolving this recursion into a loop leads to Algorithm 6, in which the lower triangle of the matrix  $B$  is overwritten with the Cholesky factor  $L$ .

By replacing matrix entries  $B(i,j)$  with  $n_b$ -by- $n_b$  blocks  $B[i,j] := B((i-1)n_b+1 : in_b, (j-1)n_b+1 : jn_b)$  we arrive at Algorithm 7 ( $N = \lceil n/n_b \rceil$  is the number of blocks).

Both algorithms require roughly  $\frac{1}{3}n^3$  flop. The operations of Algorithm 7 are done almost exclusively with level 3 BLAS.

---

**Algorithm 6** Cholesky decomposition.

---

```

for  $k = 1$  to  $n$ 
   $B(k, k) := \sqrt{B(k, k)}$ 
   $B(k + 1 : n, k) := B(k + 1 : n, k) \cdot \frac{1}{B(k, k)}$ 
   $B(k + 1 : n, k + 1 : n) := B(k + 1 : n, k + 1 : n)$ 
     $- B(k + 1 : n, k) \cdot B(k + 1 : n, k)^\top$ 

```

---



---

**Algorithm 7** Blocked Cholesky decomposition.

---

```

for  $k = 1$  to  $N$ 
   $L[k, k] :=$  Cholesky factor of  $B[k, k]$  { Algorithm 6, overwriting  $B[k, k]$  }
   $B[k + 1 : N, k] := B[k + 1 : N, k] \cdot L[k, k]^{-1}$ 
   $B[k + 1 : N, k + 1 : N] := B[k + 1 : N, k + 1 : N]$ 
     $- B[k + 1 : N, k] \cdot B[k + 1 : N, k]^\top$ 

```

---

## 5 Phase II: Methods for Tridiagonal Matrices

Once the symmetric matrix  $A$  is reduced to tridiagonal form, the eigenvalues  $\lambda_i$  and the eigenvectors  $\mathbf{v}_i$  of the tridiagonal matrix  $T$  must be found. There is a large variety of algorithms for solving this problem.

Since its invention in the 1960's, the QR iteration has been the method of choice for computing *all* eigenvalues (and, optionally, all eigenvectors). If only *selected eigenvalues* are required then bisection is the adequate method, otherwise QR-style algorithms tend to be faster. Inverse iteration may be used to compute *selected eigenvectors*, the subset being chosen after having determined and inspected the eigenvalues. In the 1990's the new divide-and-conquer algorithm, which has been developed with the aim of exposing parallelism, has proved by far superior to the QR iteration even on serial computers. It requires, however, significantly more memory than QR. Finally, there are other techniques — like homotopy algorithms — that are still in an experimental state.

Throughout this section, we will use the shorthands  $T(i, i) =: \alpha_i$  and  $T(i, i - 1) =: \beta_i$  for the diagonal and subdiagonal entries, respectively, of the symmetric tridiagonal matrix  $T$ , so

$$T = \begin{pmatrix} \alpha_1 & \beta_2 & & & \\ \beta_2 & \alpha_2 & \cdot & & \\ & \cdot & \cdot & \cdot & \\ & & \cdot & \cdot & \beta_n \\ & & & \beta_n & \alpha_n \end{pmatrix}.$$

We assume that all  $\beta_i$  are nonzero because otherwise the problem splits into two smaller subproblems that can be handled independently, cf. Section 2.5.

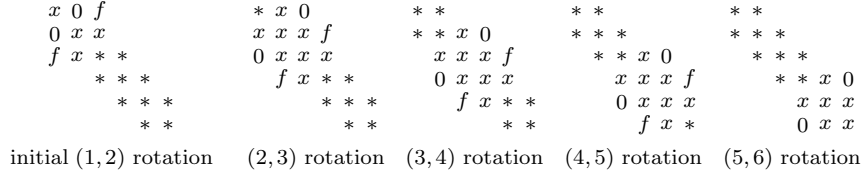


Figure 3. Rotations in one sweep of the QR iteration ( $n = 6$ ).

### 5.1 QR Iteration

The QR iteration<sup>17</sup> for symmetric tridiagonal matrices is summarized in Algorithm 8. One pass through the **repeat** loop is called a *sweep* of the iteration.

---

**Algorithm 8** QR iteration for tridiagonal matrices.

---

**repeat**

    determine a suitable shift  $\sigma \in \mathbb{R}$

    apply a rotation in the  $(1, 2)$  plane that zeroes the second entry of the vector  $(\alpha_1 - \sigma, \beta_2, 0, \dots, 0)^\top$ . This creates a fill-in element at position  $(3, 1)$ .

**for**  $i = 2$  **to**  $n - 1$

        apply a rotation in the  $(i, i + 1)$  plane that zeroes the fill-in element at position  $(i + 1, i - 1)$  and creates a new fill-in element at  $(i + 2, i)$

**until** some subdiagonal entry  $\beta_i$  becomes negligible

set  $\beta_i := 0$  and apply the whole algorithm to the submatrices  $T(1 : i, 1 : i)$  and  $T(i + 1 : n, i + 1 : n)$       { deflation }

---

Each sweep is initiated with a rotation  $T \mapsto R^\top T R$ , where the rotation angle depends on a certain parameter  $\sigma$  (the *shift*). This rotation produces a fill-in element below the subdiagonal, and a whole sequence of additional rotations is used to chase the fill element down the band until the tridiagonal structure is restored, cf. Figure 3. (See also Section 4.4 for a similar chasing strategy.)

If the shifts are chosen appropriately then the subdiagonal entries  $\beta_i$  tend to zero. Typically, the last subdiagonal entry  $\beta_n$  is the first to become negligible so that the last line in the algorithm reduces to continuing with the matrix  $T(1 : n - 1, 1 : n - 1)$ , and that one eigenvalue may be read off from the diagonal entry  $\alpha_n$ . But it is also possible that the matrix splits somewhere in the middle. The QR iteration tends to compute the eigenvalues by increasing absolute value, but this order may be broken.

The speed of convergence is determined by the shifts. For the most popular choice, *Wilkinson's shifts*, one can prove that the convergence is *global* (i.e., the subdiagonal elements are guaranteed to tend to zero) and — with rare exceptions — ultimately *cubic*. (Using Wilkinson's shifts means setting  $\sigma$  to one of the two eigenvalues of the current matrix  $T$ 's trailing 2-by-2 subblock  $\begin{pmatrix} \alpha_{n-1} & \beta_n \\ \beta_n & \alpha_n \end{pmatrix}$ , namely

the eigenvalue closer to  $\alpha_n$ .) The high order of convergence explains why — after some initial “warming up” — only two sweeps (with order-of- $n$  flop per sweep) are needed on average in order to split off another eigenvalue. Therefore only order-of- $n^2$  flop are required to compute all eigenvalues of  $T$ . If the eigenvectors of  $T$  are needed, too, then each rotation must also be applied to an  $n$ -by- $n$  matrix  $V$ ,  $V \mapsto VR$ , where  $V$  has been initialized as the orthogonal matrix  $Q_1$  that reduced the full or banded matrix  $A$  to tridiagonal form. (In the full case Algorithm 3, initialized with  $Q := I$ , can be used to build the matrix  $Q_1$  with  $\frac{4}{3}n^3$  flop.) Accumulating  $V$  requires a total of roughly  $6n^3$  flop. As the QR iteration relies completely on rotations it is backward stable.

The name of the method stems from the fact that one sweep of the iteration corresponds to first computing a *QR decomposition*  $T - \sigma I =: \tilde{Q}R$  into an orthogonal matrix  $\tilde{Q}$  and an upper triangular matrix  $R$ , and then replacing  $T$  with  $R\tilde{Q} + \sigma I$ .

This idea can also be applied to matrices with semibandwidth  $b > 1$  (indeed, even to full matrices), but for complexity reasons (one sweep then takes order-of- $b^2n$  flop without the work on the eigenvectors) such matrices are typically reduced to tridiagonal form before the iteration is started. If many eigenvalues but only a few eigenvectors of the banded matrix  $A$  are required then the following hybrid technique is may pay: First reduce  $A$  to tridiagonal form *without* accumulating the transformations, then compute the eigenvalues of  $T$ , and finally use the computed eigenvalues as “perfect” shifts in the QR iteration on the original  $A$  (with accumulation of the transformations in a matrix  $V$ , initialized as  $V = I$ ) in order to obtain the required eigenvectors.

There is a variety of mostly newer (and often more efficient) methods that are in some way similar to the tridiagonal QR iteration, in particular the LR iteration, which relies on LR (i.e., LU) decompositions of the matrix, the QL iteration, which works bottom-up instead of top-down, the Pal-Walker-Kahan QR variant for computing eigenvalues only, which requires no squares roots, and the new qd algorithms<sup>33</sup>, which essentially perform LR iteration on a factored representation of the tridiagonal matrix.

All these variants bear very limited potential for parallelism in the work on  $T$ . By contrast, the accumulation of the eigenvector matrix  $V$  is easily parallelized.

## 5.2 Bisection

Bisection is a versatile method for computing all or selected eigenvalues of  $T$ . It is based on the fact that for any  $\mu \in \mathbb{R}$ , the number  $\nu(\mu)$  of negative elements in the sequence

$$\begin{aligned} \delta_1 &:= \alpha_1 - \mu, \\ \delta_i &:= (\alpha_i - \mu) - \frac{\beta_i^2}{\delta_{i-1}}, \quad i = 2, \dots, n, \end{aligned} \tag{21}$$

is equal to the number of eigenvalues of  $T$  that are smaller than  $\mu$ . (This is just Sylvester’s law of inertia<sup>34</sup>, applied to the decomposition  $T - \mu I = LDL^\top$  with a lower triangular matrix  $L$  having all ones on the diagonal, and a diagonal matrix  $D$ .)

Therefore, given an interval  $[a, b]$  that is known to include the  $k$ -smallest eigenvalue  $\lambda_k$  of  $\mathbf{T}$ , we may locate this eigenvalue to very high precision by repeatedly cutting the interval into two halves and testing which of the halves contains  $\lambda_k$ , see Algorithm 9. Note that a suitable initial “search interval” is provided by Gershgorin’s theorem, which states that *all* eigenvalues of  $\mathbf{T}$  are contained in the interval

$$I = [\min\{\alpha_i - \rho_i : i = 1, \dots, n\}, \max\{\alpha_i + \rho_i : i = 1, \dots, n\}],$$

where  $\rho_i = |\beta_i| + |\beta_{i+1}|$ , and  $\beta_1 := 0$  and  $\beta_{n+1} := 0$  for convenience.

---

**Algorithm 9** Bisection for locating eigenvalue  $\lambda_k$  in an interval  $[a, b]$ .

---

```

compute  $\nu(a)$  and  $\nu(b)$  by building the sequences (21) for  $\mu = a$  and  $\mu = b$ 
if  $\nu(a) < k \leq \nu(b)$  { otherwise  $\lambda_k$  is not contained in  $[a, b]$  }
  while the interval width  $b - a$  is too large
    let  $c = (a + b)/2$  and compute  $\nu(c)$  by evaluating (21) for  $\mu = c$ 
    if  $\nu(c) \geq k$  {  $\lambda_k$  lies in the left half of  $[a, b]$  }
       $b := c$ 
    else {  $\lambda_k$  lies in the right half of  $[a, b]$  }
       $a := c$ 
  return  $\lambda_k \approx (a + b)/2$ 

```

---

Algorithm 9 is easily modified to compute a sequence  $\lambda_j, \dots, \lambda_k$  of consecutive eigenvalues, or all eigenvalues in a given interval  $[a, b]$ . Although bisection does not rely on orthogonal transformations (indeed, it is related to the unstable Gaussian elimination *without pivoting*), it is a perfectly stable algorithm featuring a very low backward error. In addition, bisection is efficiently parallelizable by having different processors compute disjoint subsets of the desired eigenvalues.

The major drawback of the method is its slow (namely, linear) convergence. This problem can be alleviated to some extent by using bisection only for *isolating* the eigenvalues, i.e., for narrowing down the initial interval until each subinterval contains exactly one eigenvalue. Then we switch to a superlinearly convergent root-finder, like Newton’s method<sup>25</sup> or **zeroin**<sup>8</sup>, for obtaining these eigenvalues to higher accuracy. Nevertheless, QR-type methods are superior on serial machines if more than one third, say, of the spectrum is required.

### 5.3 Inverse Iteration

Inverse iteration complements bisection in that it allows determining eigenvectors to selected previously computed eigenvalues. Inverse iteration is based on the *power iteration*, which is shown in Algorithm 10. The normalization (last line of the algorithm) is necessary to avoid overflow.

If there is a simple *dominant eigenvalue* (i.e.,  $|\lambda_{\max}| > |\lambda_i|$  for the  $n - 1$  other eigenvalues) then the vectors  $\text{sign}(\lambda_{\max}) \cdot \mathbf{v}$  tend to an eigenvector corresponding to  $\lambda_{\max}$  (the *dominant eigenvector*), and  $\|\tilde{\mathbf{v}}\|$  approaches  $|\lambda_{\max}|$ . The convergence is linear with factor  $c = |\lambda_{\max}/\lambda_{\max_2}|$ , where  $\lambda_{\max_2}$  is the eigenvalue with second largest modulus. Therefore, the more dominant  $\lambda_{\max}$  is, the faster the convergence.

Now, if  $(\lambda_i, \mathbf{v}_i)$  is an eigenpair of  $\mathbf{T}$ , and  $\mu \in \mathbb{R}$ , then  $(\lambda_i - \mu, \mathbf{v}_i)$  is an eigenpair of  $\mathbf{T} - \mu\mathbf{I}$ , and  $(1/(\lambda_i - \mu), \mathbf{v}_i)$  is an eigenpair of  $(\mathbf{T} - \mu\mathbf{I})^{-1}$ . Therefore, if  $\mu$  is a very

---

**Algorithm 10** Power iteration for approximating the dominant eigenvector.

---

select a suitable starting vector  $\mathbf{v}$

**repeat**

$$\tilde{\mathbf{v}} := \mathbf{T} \cdot \mathbf{v}$$

$$\mathbf{v} := \tilde{\mathbf{v}} / \|\tilde{\mathbf{v}}\|$$

**until**  $\mathbf{v}$  is “good enough”

---

good approximation to  $\lambda_i$  (i.e.,  $|\lambda_i - \mu| \ll |\lambda_j - \mu|$  for all  $j \neq i$ ) then  $1/(\lambda_i - \mu)$  is a strongly dominant eigenvalue of  $(\mathbf{T} - \mu\mathbf{I})^{-1}$ . This is the basis of inverse iteration, which is just the power iteration applied to  $(\mathbf{T} - \mu\mathbf{I})^{-1}$  and therefore can be obtained by substituting “solve the system  $(\mathbf{T} - \mu\mathbf{I}) \cdot \tilde{\mathbf{v}} = \mathbf{v}$ ” for the third line in Algorithm 10.

Under favorable circumstances inverse iteration is a very flexible and fast method, able to compute just the needed eigenvectors and requiring only order-of- $n$  flop per eigenvector, one order of magnitude less than QR iteration! This is the case if the eigenvalues of  $\mathbf{T}$  are well separated and if good approximations to them have been computed (e.g., with bisection or QR iteration). Then experience shows that for each eigenvalue *only one or two* iteration steps are needed to obtain an excellent approximation to the corresponding eigenvector. In addition, eigenvectors to different eigenvalues can be computed in parallel.

Problems occur with *clustered* eigenvalues, i.e.,  $\lambda_i \approx \lambda_{i+1} \approx \dots \approx \lambda_j$ . Then the orthogonality of the corresponding computed eigenvectors  $\text{fl}(\mathbf{v}_i), \dots, \text{fl}(\mathbf{v}_j)$  gets impaired because the clustered eigenvalues lead to almost identical linear systems  $(\mathbf{T} - \mu\mathbf{I})\tilde{\mathbf{v}} = \mathbf{v}$ , which in turn lead to similar solutions  $\tilde{\mathbf{v}}$ . This is particularly true if the starting vector  $\mathbf{v}$  always remains the same. Therefore, one partial remedy is to use a new random starting vector for each eigenvalue. In addition, the vectors are explicitly *orthogonalized* against each other. Suppose that we already have computed mutually orthogonal, normalized eigenvectors for the eigenvalues  $\lambda_i, \dots, \lambda_\ell$  of the cluster. Then, in each iteration step for eigenvector  $\mathbf{v}_{\ell+1}$ , we apply the *modified Gram-Schmidt process* to the intermediate vector  $\tilde{\mathbf{v}}$  before it is normalized. That is, for  $r = i, \dots, \ell$  the component of  $\tilde{\mathbf{v}}$  in direction  $\mathbf{v}_r$  is eliminated by replacing  $\tilde{\mathbf{v}}$  with  $\tilde{\mathbf{v}} - (\mathbf{v}_r^\top \tilde{\mathbf{v}}) \cdot \mathbf{v}_r$ . Note that the explicit orthogonalization may sum up to order-of- $n^3$  flop if large clusters of eigenvalues are present. Both cures cannot preclude another type of failure, which amounts to matching the eigenvectors with wrong eigenvalues<sup>23</sup>.

Currently a new variant of inverse iteration is being developed that always gives orthonormal eigenvectors *without* explicit orthogonalization<sup>12</sup>. If these efforts are successful then the combination of bisection (or QR iteration) with the new method will make all other tridiagonal eigensolvers obsolete.

#### 5.4 Divide-and-Conquer

The divide-and-conquer algorithm<sup>10,16</sup> was developed with the aim of exploiting parallelism. As it turned out, this method typically beats the QR iteration and bisection/inverse iteration even on serial machines.

The first step in the divide-and-conquer algorithm consists of “tearing” the tridiagonal matrix into two halves via a suitable rank-1-modification. For a 6-by-6 matrix, this tearing might look like

$$\begin{aligned} \mathbf{T} = \begin{pmatrix} \alpha_1 & \beta_2 & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & \\ & \beta_3 & \alpha_3 & \beta_4 & & \\ & & \beta_4 & \alpha_4 & \beta_5 & \\ & & & \beta_5 & \alpha_5 & \beta_6 \\ & & & & \beta_6 & \alpha_6 \end{pmatrix} &= \left( \begin{array}{ccc|ccc} \alpha_1 & \beta_2 & & & & \\ \beta_2 & \alpha_2 & \beta_3 & & & \\ & \beta_3 & \alpha_3 - \beta_4 & & & \\ \hline & & & \alpha_4 - \beta_4 & \beta_5 & \\ & & & \beta_5 & \alpha_5 & \beta_6 \\ & & & & \beta_6 & \alpha_6 \end{array} \right) + \begin{pmatrix} & & & & \beta_4 & \beta_4 \\ & & & & \beta_4 & \beta_4 \\ & & & & & \\ & & & & & \\ & & & & & \\ & & & & & \end{pmatrix} \\ &=: \left( \begin{array}{c|c} \mathbf{T}_1 & \\ \hline & \mathbf{T}_2 \end{array} \right) + \rho \mathbf{w} \mathbf{w}^\top \end{aligned}$$

with  $\rho = \beta_4$  and  $\mathbf{w} = (0, 0, 1, 1, 0, 0)^\top$ . In general,  $\mathbf{T}_1$  and  $\mathbf{T}_2$  are  $m$ -by- $m$  and  $(n-m)$ -by- $(n-m)$ , resp., where  $m \approx n/2$ ,  $\rho = \beta_{m+1}$ , and  $\mathbf{w} = (0, \dots, 0, 1, 1, 0, \dots, 0)^\top$ .

Then the eigendecompositions  $\mathbf{T}_1 = \mathbf{X}_1 \Delta_1 \mathbf{X}_1^\top$  and  $\mathbf{T}_2 = \mathbf{X}_2 \Delta_2 \mathbf{X}_2^\top$  are computed by applying the divide-and-conquer algorithm recursively to the smaller matrices  $\mathbf{T}_i$  (or, if these are small enough, QR iteration is used instead). This yields

$$\begin{aligned} \mathbf{T} &= \begin{pmatrix} \mathbf{X}_1 \Delta_1 \mathbf{X}_1^\top & \\ & \mathbf{X}_2 \Delta_2 \mathbf{X}_2^\top \end{pmatrix} + \rho \mathbf{w} \mathbf{w}^\top \\ &= \begin{pmatrix} \mathbf{X}_1 & \\ & \mathbf{X}_2 \end{pmatrix} \cdot \left( \begin{pmatrix} \Delta_1 & \\ & \Delta_2 \end{pmatrix} + \rho \mathbf{z} \mathbf{z}^\top \right) \cdot \begin{pmatrix} \mathbf{X}_1^\top & \\ & \mathbf{X}_2^\top \end{pmatrix} \\ &=: \mathbf{X} \cdot (\Delta + \rho \mathbf{z} \mathbf{z}^\top) \cdot \mathbf{X}^\top, \end{aligned}$$

where

$$\mathbf{z} = \mathbf{X}^\top \cdot \mathbf{w} = \begin{pmatrix} \text{last column of } \mathbf{X}_1^\top \\ \text{first column of } \mathbf{X}_2^\top \end{pmatrix}.$$

Then the eigendecomposition

$$\Delta + \rho \mathbf{z} \mathbf{z}^\top = \mathbf{Y} \Lambda \mathbf{Y}^\top \quad (22)$$

of a rank-1-perturbed diagonal matrix must be computed (see below), and finally the eigendecomposition of  $\mathbf{T}$  can be recovered via  $\mathbf{T} = (\mathbf{X} \mathbf{Y}) \cdot \Lambda \cdot (\mathbf{X} \mathbf{Y})^\top$ .

The crucial step in the overall algorithm is the efficient and stable computation of the eigendecomposition (22), the remaining operations being expensive but trivial.

Suppose for a moment that all the entries of  $\mathbf{z} = (\zeta_1, \dots, \zeta_n)^\top$  are nonzero and that the eigenvalues (i.e., the diagonal entries) of  $\Delta$  are distinct:  $\delta_1 < \delta_2 < \dots < \delta_n$ . Then it is easy to show that the eigenvalues  $\lambda_i$  of  $\Delta + \rho \mathbf{z} \mathbf{z}^\top$  are the roots of the so-called *secular equation*

$$f(\lambda) := 1 + \rho \sum_{i=1}^n \frac{\zeta_i^2}{\delta_i - \lambda} = 0, \quad (23)$$

and that for any eigenvalue  $\lambda_i$  of  $\Delta + \rho \mathbf{z} \mathbf{z}^\top$ ,

$$\mathbf{y}_i := (\Delta - \lambda_i \mathbf{I})^{-1} \cdot \mathbf{z} \quad (24)$$

is a corresponding eigenvector. (Note that computing  $\mathbf{y}_i$  involves just an appropriate entry-wise scaling of  $\mathbf{z}$ .)

Eq. (23) implies that the  $\delta_i$  “interlace” the sought eigenvalues, i.e., for  $\rho > 0$  we have

$$\delta_1 < \lambda_1 < \delta_2 < \lambda_2 < \delta_3 < \dots < \delta_n < \lambda_n .$$

This property, together with the fact that on each interval  $(\delta_i, \delta_{i+1})$  the function  $f$  can be approximated by simple rational expressions, leads to a globally and quadratically convergent root finder that can compute all the  $\lambda_i$  in order-of- $n^2$  time.

Unfortunately, Eq. (24) is not an adequate means to compute the eigenvectors because orthogonality is severely impaired in the presence of close  $\delta_i$ . It took more than ten years from the invention of the divide-and-conquer method until a technique was discovered<sup>19</sup> that did not need resorting to extended precision in the eigenvector computations. Roughly speaking, we can compute orthogonal eigenvectors for  $\Delta + \rho \mathbf{z} \mathbf{z}^\top$  by applying Formula (24) to a slightly modified vector  $\mathbf{z}$ .

In practice our assumption that all the  $\zeta_i$  are nonzero and all the  $\delta_i$  are distinct is seldom fulfilled. The superiority of the divide-and-conquer method comes from the fact that it even can take advantage from a violation of this assumption. In fact, if some  $\zeta_i$  is zero then  $\delta_i$  and the  $i$ th column of  $\mathbf{X}$  already are an eigenpair of the tridiagonal matrix  $\mathbf{T}$ . Similarly, if some of the  $\delta_i$  are (almost) identical then all but one of them are also very good approximations to eigenvalues  $\lambda_i$  of  $\mathbf{T}$ , and the corresponding eigenvalues can be computed cheaply. Thus, the eigenvalue problem for  $\Delta + \rho \mathbf{z} \mathbf{z}^\top$  and the ensuing multiplication  $\mathbf{X} \cdot \mathbf{Y}$  are effectively reduced in size. Fortunately this type of *deflation* is quite frequent.

Parallelism can be exploited in two ways. First, all the solutions of the secular equation and the corresponding eigenvectors of  $\Delta + \rho \mathbf{z} \mathbf{z}^\top$  may be computed independently from each other, and second, the matrix–matrix product  $\mathbf{X} \cdot \mathbf{Y}$  lends itself naturally to a coarse-grained parallelization.

There are also attempts to apply the divide-and-conquer technique to banded matrices<sup>3</sup>, but the resulting algorithms are still highly experimental.

### 5.5 Homotopy Methods

Homotopy methods<sup>29</sup> try to follow the paths of the eigenvalues and eigenvectors through a whole sequence of matrices

$$\mathbf{T}_0 \mapsto \mathbf{T}_1 = \mathbf{T}_0 + \theta_1(\mathbf{T} - \mathbf{T}_0) \mapsto \dots \mapsto \mathbf{T}_{k-1} = \mathbf{T}_0 + \theta_{k-1}(\mathbf{T} - \mathbf{T}_0) \mapsto \mathbf{T}_k = \mathbf{T} ,$$

where  $\mathbf{T}_0$  is some initial matrix whose eigensystem is readily computed (e.g., a diagonal matrix),  $\mathbf{T}$  is the matrix whose eigensystem is sought, and  $0 = \theta_0 < \theta_1 < \dots < \theta_{k-1} < \theta_k = 1$ . If the step-sizes  $\theta_i - \theta_{i-1}$  are small enough then the eigenvalues of  $\mathbf{T}_{i-1}$  are good starting values for computing the eigenvalues of  $\mathbf{T}_i$ , and inverse iteration for  $\mathbf{T}_i$ ’s eigenvectors can be started with the eigenvectors of  $\mathbf{T}_{i-1}$ .

The same idea may also be applied to other (e.g., banded) matrices and to the generalized eigenvalue problem. Note that the homotopy methods are still in an experimental state.



## 6 Methods Without Initial Tridiagonalization

In contrast to the methods described above, the algorithms discussed in this section do not rely on an initial tridiagonalization of the matrix, but apply the iterative process to the matrix  $A$  itself.

### 6.1 Jacobi's Method

Jacobi's method is based on the idea of reducing  $A$ 's "off-diagonal norm",

$$\text{off}(A) := \sqrt{\sum_{i=1}^n \sum_{j=1, j \neq i}^n A(i, j)^2},$$

until it is negligibly small. Then we can read off  $A$ 's eigenvalue from the diagonal entries and the eigenvectors from the columns of the orthogonal transformation matrix that was used to attain the almost-diagonal form.

This form is achieved by repeatedly zeroing selected entries  $A(i, j)$  with suitable rotations. If we choose the rotation angle  $\theta$  such that

$$\tan \theta = \frac{\text{sign}(\tau)}{|\tau| + \sqrt{1 + \tau^2}}, \quad \text{where } \tau = \frac{A(i, i) - A(j, j)}{2A(i, j)},$$

then a short computations shows that the two-sided rotation  $A \mapsto R^\top A R$  with  $R = R(i, j, \theta)$  indeed zeroes  $A$ 's  $(i, j)$  and  $(j, i)$  entries and that  $\text{off}(A)$  drops by  $2 \cdot A(i, j)^2$ . Note that again the rotation angle is not needed explicitly because the parameters  $c$  and  $s$  may be obtained via

$$c = \frac{1}{\sqrt{1 + t^2}}, \quad s = t \cdot c.$$

Unfortunately, zeros introduced this way do not persist but are made nonzero again in later rotations. Thus most entries of  $A$  have to be made zero several times during the whole process. There are many different strategies for selecting the order of the entries  $(i, j)$  to be zeroed.

Obviously, zeroing the off-diagonal entry  $A(i, j)$  with the largest absolute value will lead to the largest reduction of  $\text{off}(A)$ . This is the *classical Jacobi method*<sup>22</sup>, which is slowed down by organizational overhead since the roughly  $n^2/2$  comparisons for determining the maximum entry cost by far more time than the  $6n$  ensuing arithmetic operations.

Therefore most often cheaper schemes are used, in particular the *row cyclic* and *column cyclic* elimination orders. In the former, the entries of the strictly lower triangle of  $A$  are made zero row-by-row, that is, in the order  $A(2, 1), A(3, 1), A(3, 2), A(4, 1), A(4, 2), A(4, 3), \dots, A(n, 1), A(n, 2), \dots, A(n, n-1)$ . When all these entries have been made zero once (this is called a *sweep* of the method), then the process is started anew.

It is easy to see that each rotation in the classical method must reduce  $\text{off}(A)$  by a factor  $\leq 1 - \frac{2}{n(n-1)} < 1$ , thus implying at least linear convergence. A more involved analysis reveals that the convergence is indeed much faster, namely

quadratic<sup>34</sup>. This is also true for the cyclic schemes, provided that some precautions concerning the rotation angles are taken.

Experimental evidence suggests that in practice roughly  $\log n$  sweeps are necessary to achieve adequate accuracy. As the methods discussed in Sections 4 and 5 require much less work (corresponding to just two Jacobi sweeps), Jacobi's method is usually not competitive. If, however, the matrix  $A$  is already strongly *diagonally dominant* (i.e., its diagonal entries are much larger than the off-diagonals) then Jacobi's method needs only a few sweeps to converge and may even beat the reduction-based methods.

Two other facts have revived the interest in Jacobi's method. First, it was observed that in some cases this algorithm — carefully implemented — can deliver much more accurate eigensystems than the other techniques. And second, appropriate cyclic elimination schemes allow exploiting parallelism by applying several rotations simultaneously<sup>31</sup>.

There are also variants of Jacobi's method for the generalized eigenvalue problem, but their convergence properties are not sufficiently known.

## 6.2 Invariant Subspace Decomposition

The *invariant subspace decomposition algorithm*<sup>21</sup> (ISDA) is in some sense dual to the bisection/inverse iteration approach since it extracts information about the eigenvectors *before* the eigenvalues. This algorithm relies on the fact that eigenvectors are invariant under polynomial transformations of the matrix: If  $(\lambda, \mathbf{q})$  is an eigenpair of  $A$  then for any polynomial  $p$ ,  $(p(\lambda), \mathbf{q})$  is an eigenpair of  $p(A)$ .

First, lower and upper bounds for  $\text{spec}(A)$  are determined (e.g., with Gershgorin's theorem for full matrices), and then a linear transformation  $A \mapsto \alpha A + \beta I =: A_1$  is applied that maps the lower half of the eigenvalues,  $\lambda_1, \dots, \lambda_{n/2}$ , into the interval  $[0, \frac{1}{2}]$  and the upper half of the eigenvalues into  $[\frac{1}{2}, 1]$ .

Then further polynomial transformations  $A_k \mapsto p_k(A_k) =: A_{k+1}$  are applied, where the polynomial  $p_k$  is designed such that  $p_k(x) \approx 0$  for  $x \in [0, \frac{1}{2})$  and  $p_k(x) \approx 1$  for  $x \in (\frac{1}{2}, 1]$ , i.e.,  $p_k$  pushes the lower half of the eigenvalues toward 0 and the upper half toward 1. This process is repeated until  $A_k \approx A_{k-1}$ , implying that all eigenvalues of the final  $A_k$  are approximately 0 or 1.

Next, a so-called *rank-revealing QR decomposition*  $A_k = QR$  of this matrix into an orthogonal matrix  $Q$  and an upper triangular matrix  $R$  is determined. Then the first (last)  $n/2$  columns of  $Q$  are orthonormal eigenvectors to the eigenvalue 1 (0, respectively) of  $A_k$ . This implies

$$Q^T A Q = \begin{pmatrix} A' & \\ & A'' \end{pmatrix},$$

where  $A'$  and  $A''$  are symmetric  $(n/2)$ -by- $(n/2)$  matrices. According to Section 2.5, the eigenvalues and eigenvectors of  $A$  may now be obtained from the eigendecompositions of  $A'$  and  $A''$ , which in turn are computed by recursively applying the ISDA to these two matrices, or with QR iteration if  $A'$  and  $A''$  are small enough.

This approach does most of its computations with matrix-matrix products (e.g., in the evaluation of  $p_k(A_k)$ ) and therefore achieves high Mflop/s rates. In addition,

efficient parallelization is possible. On the other hand, the ISDA requires several times more **flop** than the reduction-based techniques and is therefore not competitive on serial machines.

## 7 Available Software

Whenever possible, eigenvalue solvers from established libraries should be used because much effort went into optimizing their performance and making them robust (e.g., appropriate scalings for “balancing” the eigenvalues). There are many pitfalls awaiting the ambitious but unexperienced programmer.

### 7.1 Serial and Shared-Memory Machines

As modern software for dense or banded matrices relies heavily on the BLAS, obtaining an optimized implementation of the BLAS is extremely important. For most high-performance machines, optimized BLAS are provided by the manufacturer. If this is not the case, the public-domain BLAS from the ATLAS project (<http://www.netlib.org/atlas>) are a viable alternative. They often perform almost as well as (and sometimes even better than) proprietary implementations.

The public-domain *LAPACK* library<sup>1</sup> (<http://www.netlib.org/lapack>) contains optimized implementations for almost all non-experimental algorithms described in this article. (LAPACK also solves linear systems and least squares and related problems, and most routines are available for real and complex matrices in either single-precision or double-precision arithmetic.) Preferably the library’s comfortable *driver routines* should be used. E.g., *DSYEV* computes the eigendecomposition of a full matrix, whereas *DSBEV* handles the banded case. Similar drivers exist for the generalized (full or banded) eigenvalue problem, as well as so-called “expert” drivers with additional options (e.g., computing only parts of the eigendecomposition). In addition, the computational routines may be called directly. Thus, *DSYTRD* (*DSBTRD*) reduces a full (banded) matrix to tridiagonal form with Algorithm 2 (with a modified version of Algorithm 4).

At <http://www-unix.mcs.anl.gov/prism>, additional software for the ISDA (Section 6.2) and for the two-phase reduction (mentioned at the end of Section 4.2) is available.

Because of the high quality of the LAPACK library, many vendor-supplied numerical packages and commercial libraries (like NAG) are based on these routines. This is also the case for *multithreaded* parallel libraries that come with shared-memory parallel machines. Here the parallelism is often confined within the BLAS.

### 7.2 Distributed-Memory Systems

For distributed-memory machines with the message-passing programming model, the situation is more complicated. The direct analogue to LAPACK is the *ScaLAPACK* library<sup>7</sup> (<http://www.netlib.org/scalapack>), which contains the functionality of many LAPACK routines. ScaLAPACK is based on the *PBLAS*, a parallelized implementation of the BLAS. ScaLAPACK contains solvers for the symmetric standard and generalized eigenvalue problem with full matrices, but *no*

*banded solvers*. Another problem with this library is the fact that inverse iteration performs explicit orthogonalization only against eigenvectors within the same processor. Thus orthogonality may be lost. (ScaLAPACK also includes QR iteration, which does not suffer from this problem.)

*PeIGS* (<http://www.emsl.pnl.gov:2080/docs/global/peigs.html>) is another library offering eigensystem functionality comparable to that of ScaLAPACK; in particular, banded problems are not addressed. While ScaLAPACK relies on a two-dimensional data layout (i.e., the matrices are split along rows *and* columns), PeIGS works with matrices that are distributed by whole columns or whole rows. Such distributions are typically inferior. On the other hand, inverse iteration in PeIGS orthogonalizes against all eigenvectors of a cluster and therefore gives better results.

Finally, there are packages that do not contain complete eigensolvers but rather provide infrastructure for easily putting such methods together, namely *PLAPACK*<sup>38</sup> (<http://www.cs.utexas.edu/users/plapack>) and *Global Arrays* (<http://www.emsl.pnl.gov:2080/docs/global/ga.html>). Both packages facilitate the manipulation of distributed matrices.

## References

1. E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney and D. Sorensen, *LAPACK Users' Guide*, 3rd ed., SIAM, Philadelphia, PA (1999).
2. ANSI/IEEE, IEEE Standard for Binary Floating Point Arithmetic, Std. 754, New York (1985).
3. P. Arbenz, Divide and conquer algorithms for the bandsymmetric eigenvalue problem, *Parallel Comput.* **18**, 1105–1128 (1992).
4. J. Barlow and J. Demmel, Computing accurate eigensystems of scaled diagonally dominant matrices, *SIAM J. Numer. Anal.* **27(3)**, 762–791 (1990).
5. C. Bischof, B. Lang, and X. Sun, Parallel tridiagonalization through two-step band reduction, in *Proceedings of the Scalable High-Performance Computing Conference, Knoxville, Tennessee, May 23–25, 1994*, IEEE Computer Society, Los Alamitos, CA, 23–27 (1994).
6. C.H. Bischof and C.F. Van Loan, The WY representation for products of Householder matrices, *SIAM J. Sci. Stat. Comput.* **8 (1)**, s2–s13 (1987).
7. L.S. Blackford, J. Choi, A. Cleary, E. d'Azevedo, J. Demmel, I. Dhillon, J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker and R.R. Whaley, *ScaLAPACK Users' Guide*, SIAM, Philadelphia, PA (1997).
8. R.P. Brent, *Algorithms for Minimization Without Derivatives*, Prentice-Hall, Englewood Cliffs, NJ (1973).
9. C.R. Crawford, Reduction of a band symmetric generalized eigenvalue problem, *Comm. Assoc. Comp. Mach.* **16**, 41–44 (1973).
10. J.J.M. Cuppen, A divide and conquer method for the symmetric tridiagonal eigenproblem, *Numer. Math.* **36**, 177–195 (1981).
11. J.W. Demmel, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, PA (1997).

12. I.S. Dhillon, *A New  $O(n^2)$  Algorithm for the Symmetric Tridiagonal Eigenvalue/Eigenvector Problem*, Ph.D. thesis, Computer Science Division, University of California at Berkeley (1997).
13. J.J. Dongarra, J. Du Croz, S. Hammarling, and I. Duff, A set of level 3 basic linear algebra subprograms, *ACM Trans. Math. Softw.* **16** (1), 1–17 (1990).
14. J.J. Dongarra, J. Du Croz, S. Hammarling, and R.J. Hanson, An extended set of FORTRAN basic linear algebra subprograms, *ACM Trans. Math. Softw.* **14** (1), 1–17 (1988).
15. J.J. Dongarra, S.J. Hammarling, and D.C. Sorensen, Block reduction to condensed forms for eigenvalue computations, *J. Comput. Appl. Math.* **27**, 215–227 (1989).
16. J.J. Dongarra and D.C. Sorensen, A fully parallel algorithm for the symmetric eigenvalue problem, *SIAM J. Sci. Stat. Comput.* **8** (2), s139–s154 (1987).
17. J.G.F. Francis, The QR transformation: A unitary analogue to the LR transformation, part I and II, *Computer J.* **4**, 265–272 and 332–345 (1961/62).
18. G.H. Golub and C.F. Van Loan, *Matrix Computations*, 3rd ed., The Johns Hopkins University Press, Baltimore, MD (1996).
19. M. Gu and S.C. Eisenstat, A stable and efficient algorithm for the rank-1 modification of the symmetric eigenvalue problem, *SIAM J. Matrix Anal. Appl.* **15**, 1266–1276 (1994).
20. N.J. Higham, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, PA (1996).
21. S. Huss-Lederman, A. Tsao and G. Zhang, A parallel implementation of the invariant subspace decomposition algorithm for dense symmetric matrices, in *Proc. Sixth SIAM Conf. on Parallel Processing for Scientific Computing in Norfolk, VA*, SIAM, Philadelphia, PA (1993).
22. C.G.J. Jacobi, Über ein leichtes Verfahren, die in der Theorie der Säculärstörungen vorkommenden Gleichungen zu lösen, *Crelle's J.* **30**, 51–94 (1846).
23. E.R. Jessup and I.C.F. Ipsen, Improving the accuracy of inverse iteration, *SIAM J. Sci. Stat. Comput.* **13**, 550–572 (1992).
24. L. Kaufman, Banded eigenvalue solvers on vector machines, *ACM Trans. Math. Softw.*, **10** (1), 73–86 (1984).
25. D. Kincaid and W. Cheney, *Numerical Analysis*, 2nd ed., Brooks/Cole Publishing Company, Pacific Grove, CA (1996).
26. B. Lang, A parallel algorithm for reducing symmetric banded matrices to tridiagonal form, *SIAM J. Sci. Comput.* **14** (6), 1320–1338 (1993).
27. B. Lang, Using level 3 BLAS in rotation-based algorithms, *SIAM J. Sci. Comput.* **19** (2), 626–634 (1998).
28. C.L. Lawson, R.J. Hanson, D.R. Kincaid, and F.T. Krogh, Basic linear algebra subprograms for FORTRAN usage, *ACM Trans. Math. Softw.* **5** (3), 308–323 (1979).
29. T.Y. Li and N. Rhee, Homotopy algorithm for symmetric eigenvalue problems, *Numer. Math.* **55**, 265–280 (1989).
30. The MathWorks, Inc., *MATLAB Reference Guide*, Natick, MA (1999).
31. J.J. Modi, *Parallel Algorithms and Matrix Computation*, Clarendon Press, Oxford (1988).

32. K. Murata and K. Horikoshi, A new method for the tridiagonalization of the symmetric band matrix, *Information Processing in Japan* **15**, 108–112 (1975).
33. B.N. Parlett, The new qd algorithms, *Acta Numerica*, 459–491 (1995).
34. B.N. Parlett, *The Symmetric Eigenvalue Problem*, SIAM, Philadelphia, PA (1998). Updated reprint of the 1980 Prentice-Hall edition.
35. W. Rath, Fast Givens rotations for orthogonal similarity transformations, *Numer. Math.* **40**, 47–56 (1982).
36. R. Schreiber and C. Van Loan, A storage efficient WY representation for products of Householder transformations, *SIAM J. Sci. Stat. Comput.* **10** (1), 53–57 (1989).
37. H.R. Schwarz, Tridiagonalization of a symmetric band matrix, *Numer. Math.* **12**, 231–241 (1968).
38. R.A. van de Geijn, *Using LAPACK*, MIT Press, Cambridge (1997).
39. J. Wilkinson, *Rounding Errors in Algebraic Processes*, Prentice-Hall (1964).
40. J.H. Wilkinson, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford (1965).